MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS 1963-A

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br>AFIT/GOR/OS/81D-2 | 2. GOVT ACCESSION NO.<br>AD-A111 428 | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br>A MONTE CARLO STUDY OF DIMENSIONALITY ASSESSMENT AND FACTOR INTERPRETATION IN PRINCIPLE COMPONENT ANALYSIS | | 5. TYPE OF REPORT & PERIOD COVERED<br>MS Thesis |
|  | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s)<br>Kenneth W. Bauer Jr.<br>Captain, USAF | | 8. CONTRACT OR GRANT NUMBER(s) |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS |
| 11. CONTROLLING OFFICE NAME AND ADDRESS | | 12. REPORT DATE<br>December 1981 |
|  | | 13. NUMBER OF PAGES<br>72 |
| 13. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br>UNCLASSIFIED |
|  | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

FEB 2 2 1982

16. DISTRIBUTION STATEMENT (of this Report)

Approved for public release: distribution unlimited

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)

**28 JAN 1982**

Air Force Institute of Technology (ATC)
Wright-Patterson AFB, OH 45433

18. SUPPLEMENTARY NOTES

APPROVED FOR PUBLIC RELEASE AFR 190-17.

FREDRIC C. LYNCH, Major, USAF
Director of Public Affairs

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

Principle component analysis
Factor analysis
Monte Carlo study

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

This study addresses two steps in a process known as principal components analysis using Monte Carlo techniques. An analysis is presented of two popular dimensionality assessment techniques, Kaiser's criterion and Catell's scree test. The factor interpretation issue is addressed through a regression study in which the grand mean square error between population and sample

factor loading matrices is predicted.  The notion of a complexity index is also introduced.

A MONTE CARLO STUDY OF

DIMENSIONALITY ASSESSMENT AND

FACTOR INTERPRETATION IN

PRINCIPLE COMPONENT ANALYSIS

AFIT/GOR/OS/81D-2      Kenneth Bauer
                       Captain   USAF

82 02 18 112

Thesis

A MONTE CARLO STUDY OF DIMENSIONALITY
ASSESSMENT AND FACTOR INTERPRETATION
IN PRINCIPLE COMPONENT ANALYSIS

by

Kenneth Bauer
Captain   USAF

Prepared in partial
fulfillment of the
requirements for a
Masters Degree

December 1981

School of Engineering
Air Force Institute of Technology
Wright-Patterson Air Force Base
Ohio

## Preface

I wish to thank both my advisor, Lt. Col. Charles McNichols, and my reader, Dr. B. N. Nagarsenkar, for their patience and for allowing me to be creative.

I would also like thank Capt. Don Turos who rescued me at the last moment when the final production of this thesis seemed doubtful. Grateful thanks are also extended to Capt. Mike Cox, who taught me virtually everything I know about the CDC 6000 computer.

I would like to dedicate this effort to my Dad, who would have been quite proud and my new son Scott, who helped out by being born healthy.

## Abstract

This study addresses two steps in  a process known as principle components analysis using Monte Carlo techniques. An analysis is presented of two popular dimensionality assessment techniques, Kaiser's criterion and Catell's Scree test. The factor interpretation issue is addressed through a regression study in which the grand mean square error between population and sample factor loading matrices is predicted. The notion of a complexity index is also introduced.

# TABLE OF CONTENTS

# List of figures

## Introduction

### Background

Factor analysis is a widely used multivariate data analysis technique. This technique allows an analyst to investigate the underlying structure of a set of variables over which data has been gathered.

Factor analysis has seen its most concentrated application in the behavioral sciences. Even though the methods and models of factor analysis are of a statistical nature, factor analysis was developed mainly by psychologists (Joreskog,1979) and as such a literature search in the area of factor analysis will lead one to such journals as Psychometrika and Psychological Reports.

Objectives of Factor Analysis. One object of factor analysis is to determine the underlying dimensionality of a process by finding independent factors which are highly related to one or more of the variables in question. The process by which an investigator determines (via factor analysis) the dimensionality of a set of data will be called, for purposes of this report, the dimensionality assessment. After the dimensionality assessment has been made it would be desirable to be able to give a simple interpretation to each of the factors (McNichols,1980). At this point a process know as rotation is used in an

1

attempt to develop a simpler more easily interpretable
solution. So, another objective for the investigator is
to interpret the extracted factors correctly. This
objective will be called, again for the purposes of
this report, factor interpretation.

Principal Component Analysis. This report will
deal exclusively with a methodology known as principal
components analysis (PCA). Psychologists draw the
following distinctions between factor analysis and
principal components analysis. In factor analysis an
attempt is made to find a certain number of factors,
fewer than the number of variables, such that the
intercorrelations between the variables is reproduced
exactly. In PCA ,independent factors are extracted from
the data until a sufficient proportion of the total
variance exhibited by the data is reproduced by the
factors. Hence, PCA is said to be variance oriented
while factor analysis is said to be correlation
oriented (Joreskog,1979). For purposes of this report
PCA will be considered a "factor analytic" procedure.

Computational Procedure. The PCA computational
procedure consists of three steps (1) the preparation
of the correlation matrix, (2) the extraction of the
initial factors - from whence a dimensionality
assessment is made (this step is where the investigator
explores the possibility of data reduction) and (3)
rotation to a term... ac .tion - a search for simple

2

and interpretable factors (Nie,1975) (McNichols,1980).

1. Preparation of the correlaton matrix. To prepare the correlation matrix the following procedure is used. First let n obsevations of some random variable be given by a nx1 vector $\underline{x}_1$.

$$\underline{x}_1 = \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ \vdots \\ x_{n1} \end{pmatrix}$$

Thus we have n observations on a single variable. If we have k such variables we can assemble them in a matrix such that each column represents n observations on the k variables. This data matrix, X, is a nxk matrix.

$$X = \begin{bmatrix} x_{11} & \cdots\cdots\cdots & x_{1n} \\ x_{21} & & \vdots \\ \vdots & & \vdots \\ \vdots & & \vdots \\ x_{n1} & \cdots\cdots\cdots & x_{nk} \end{bmatrix}$$

(Where the subscript i denotes the ith observation and

3

the subscript j denotes the jth variable.)

if each $x_{ij}$ is standardized ie. set equal to

$$\frac{x_{ij} - \overline{x}_j}{s_j}$$

(where $\overline{x}_j$ is the sample mean of the jth variable and $s_j$ is sample standard deviation of the jth variable) and if the matrix multiplication $X^T X$ is performed then the resultant product is the sample correlation matrix.

2. Extraction of the initial factors. The second step in the process is to extract the eigenvalues from $X^T X$. The extracted eigenvalues, when rank ordered, are the basis for most of the currently popular dimensionality assessment techniques. Once this assessment has been made the eigenvectors associated with the retained eigenvalues are normalized. these vectors are then multiplied by the square roots of their respective eigenvalues and assembled in matrix form. This matrix is called the factor loadings matrix, this is because each coefficient can be said to represent a "loading" of that particular variable to the factor.

3. Rotation to terminal solution. Since the normalized eigenvectors (or principal components as they are sometimes called) are by definition mutually orthogonal, the PCA methodology has,therefore, yields a

4

basis for the data space whose dimension is less than
or equal to the rank of the correlation matrix. The
loadings matrix , or factor structure matrix, is not
unique. That is to say, the factor structure matrix may
be rotated arbitrarily and will still explain the same
amount of total variance. Rotation schemes have been
developed to "clean up" the factor structure matrix and
hence aid in factor interpretation (McNichols,1980).
This rotation for interpretation is the third and final
step.

For a more detailed mathematical and statistical
development see (Lawley and Maxwell,1971),
(Harman,1967), (Harris,1975).

## Problem Statement

One major problem facing an investigator is trying to determine how many factors to retain. Unfortunately, the dimensionality assessment procedure for PCA has not been rigidly defined. Several alternatives are available to an investigator. These include the Scree test (catell, 1966), Horn's test (Horn, 1965), and the "fraction of variance explained" test (Kaiser, 1960). There does not appear to be evidence that any one of these is superior to the others. It would be of interest to determine which test is the most powerful.

Recent research suggests that factor analytic procedures can be significantly influenced by sample size, the number of variables, number of inherent factors, complexity of the inherent factor structure, and the interactions between these. The following questions are raised:

1)   To what extent is the dimensionality assessment biased by these influences? How is the variance of the dimensionality assessment affected?

2)   How is factor interpretation affected by these influences? More specifically, how is the size of the mean square error (MSE) (or perhaps RMS error) of sample loadings affected? Can MSE (or RMS errors) be reasonably predicted as a function of the forementioned influences?

6

3) An investigator beginning a factor analytic study starts with limited information. Assuming he knows the correlation matrix of the sample data; is there some rule of thumb based on the condition number of this matrix, the number of variables, and the sample size which might aid him in estimating errors he might expect in a sample factor loadings matrix?

## Review of the literature

Browne (1968) states,"There is no statistical test of significance for the number of factors applicable to the principal factor estimates and rules of thumb are generally used for estimating the number of factors."

It is important to note that factor analysis is not devoid of statistical inference tools. Cliff and Hamburger (1967) assert that

> " The results available from  statistical theory, while useful, leave a large area where the needs of the investigator are unsatisfied.  The statistical tests for the number of factors are all tied, naturally enough, to the repective methods for estimating factors. Moreover , these methods of estimation are either computationally arduous .... or unfamiliar .... Consequently, they are rarely used and so the statistical tests are rarely applied. More important than this is the fact that the number of factors in a given matrix is only one of the many concerns of the investigator. He is interested in a wide variety of statistical questions, and he is interested in them as they arise in the methods of factor analysis currently in use.

> Factor analysis is not the only field of statistics where questions of practical interest have been too complex or too difficult to specify mathematically for analytic solution. In such instances it is fairly common practice to take the Monte Carlo approach, in which samples from some specified population are generated by some random process."

Tucker (1964) and Linn (1968) used a Monte Carlo approach to study what Psychologists term pyschometric error. Psychometric error arises when, say, a given battery of tests are measuring factors different from

8

those factors it was designed to measure.

Tucker's approach entailed the perturbation of idealized common factor loadings. Tucker constructed what he called a "formal model". In the formal model, the idealized loadings were changed by random deviations. Small loadings on over 100 additional factors were also generated by a random process. In this test, then, the common factor loadings were perturbed and this effect was compounded by the presence of a large number of nuisance factors. This expanded factor loading matrix was then multiplied by its transpose to yield a correlation matrix. This correlation matrix was then factored to yield a factor loading matrix. This, then, is the factor loading matrix of what Tucker calls the "simulation model". Tucker asserts that by comparing the factor loading matrix of the "formal model" to that of the "simulation model" one can measure the degree of psychometric error between the idealized factor structure and the perturbated factor structure.

Horn (1965) suggested a novel application of the Monte Carlo procedure to aid in dimensionality assessment. Horn's procedure is as follows:

1) Given a data matrix of m measurements on n variables, the sample correlation matrix is prepared and its eigenvalues are extracted.

2) Generate n samples of m independent

9

observations from a normally distributed population of random numbers. Prepare this correlation matrix and extract its eigenvalues. This process is to be repeated K times to yield K sets of eigenvalues. Each set of eigenvalues is rank ordered and averages of the largest eigenvalue, second largest eigenvalue, etc. are computed across the K sets. (Note that K should be chosen large enough to ensure a good estimate for the means).

3) Compare, in rank order, the eigenvalues due to the real world data with the eigenvalues of the K randomly generated correlation matrices. Pick the smallest eigenvalue from the real world data correlation matrix whose value is larger than its counterpart from the randomly generated matrices.

Horn argues that this procedure will reduce the number of factors that would have been retained by, say, the Guttman (1954) weaker lower bound. (Guttman's weaker lower bound is more commonly referred to as Kaiser's criterion. Kaiser (1960) first adapted a procedure of retaining all factors whose eigenvalues were all greater than or equal to one.). Horn feels that too many factors will have been retained by the other criteria due to the fact that these tests ignore sampling error and the error which Horn refers to as least squares "capitalization". Least squares "capitalization" refers to the fact that, in factor

10

analysis, the first derived factor is constructed in such a fashion as to take up as much of the variance present as possible (in the least squares sense) and as such "capitalizes" on the chance fluctuations in a particular set of data. Horn, then (to use engineering vernacular), claims that the "real " eigenvalues ride above those eigenvalues due to the inherent "noise" of the process being investigated.      Horn makes no pretensions, however , as to the validity of his procedure and essentially presents his rationale as a rule of "thumb". Although Horn's reasoning appealed to this author, experts in the field have not received the rationale with open arms. Cliff and Hamburger (1967) go so far as to state that Horn's arguments are purely "verbal" and present a hypothetical counter example. The counter example goes as follows: suppose the experimental situation is such that there is a very large common factor present and a second much smaller one. The eigenvalue of the large factor may be so large that all succeeding eigenvalues are much less in magnitude than one. In the random data we should find several eigenvalues greater than one (approximately 1/2 of them). Hence, when the eigenvalues comparisions are made we will only retain the single large factor. So, in this counter example, Horn's rationale has underestimated the number of true factors. In defense of Horn, however, one might wonder how important that

11

other small factor was to the analysis, that is to say, what is the penalty of ignoring such small factors?

Browne (1968) published a lengthy Monte Carlo study in which he examined the effects that increasing sample size, and increasing the number of variables while holding constant the number of factors had on estimates of factor loadings. Browne compared and contrasted several factor analytic techniques (PCA was not analyzed, but a hybrid technique which placed the communalities on the diagonal of the correlation matrix was examined). The same population factor matrix was used throughout his study. This matrix had 16 variables which loaded on 4 factors with communalities ranging from .9 to .1. Using a method due to Odel and Feiveson (1966), random correlation matrices were generated and factor analyzed using five different techniques. Browne presents extensive tables which show various attributes of the sample loadings he obtained. Browne demonstrates the superiority of of a technique called "maximum likelihood" to estimate factor loadings. He also studied the dimensionality assessment problem and found that although none of the methods proved completely satisfactory, the decision procedure based on a sequence of likelihood ratio tests and the criterion of number of eigenvalues greater than one of the sample correlation matrix gave the best results of the methods considered. The other criterion considered was due to

12

Saunders (1960), this technique involves taking, at each iteration of a factor analytic technique called Thomson's method (Thomson,1934), the number of positive eigenvalues greater than the absolute value of the smallest eigenvalue as a criterion for the number of factors to be used for the following iteration.

Linn (1968) suggested an approach to the dimensionality assessment problem which is similar in spirit to Horn's procedure. Instead of analyzing the real variables and the randomly generated variables separately,they are analyzed jointly. In Linn's procedure k new variables are introduced to the real data and if there are m observations on each of the original variables, m independent observations are generated for the k new variables. All the observations across both sets of the variables are then used to prepare a correlation matrix which is subsequently factor analyzed. The dimensionality assessment is then based on those factors which are primarily composed of real variables. The factors, whose main loadings are those of the generated variables, are to be discarded, then, as random factors. Linn (1968) expanded the scope of his original 1964 study. Linn's 1964 study was limited by the fact that only two observed matrices were used and each of these was based on a sample of size 80.

Catell (1966) suggests a brief, easy to apply test

13

for dimensionality assessment. The test entails simply graphing the eigenvalue magnitudes versus the factor number (the factor number, for instance, of the factor corresponding to the largest eigenvalue is 1; the second largest eigenvalue is 2, etc.). Catell noticed that often in empirical data a "break" was exhibited in the eigenvalue curves. This break, Catell reasoned, signalled the the beginning of the trivial factors and started a more gradually sloping linear trend in the eigenvalue curve which resembled a scree line. A scree line is the shape that rocks sliding off a hill will assume at the bottom of the hill. Although this test is not Monte Carlo in nature it is mentioned because Horn (1965) and Linn (1968) both noticed this break in eigenvalues magnitude.

Hamburger (Cliff and Hamburger,1967) performed a limited Monte Carlo study to examine the apparent break in eigenvalue curves. To be more specific, Hamburger sought to find a break where two adjacently ranked eigenvalues are sharply different in size while on both sides of the break the decrease is more gradual. He reports that when sample sizes as are high as 400, then using the break in eigenvalue magnitude as a decision rule for factor retention was flawless (at least for the matrices examined in the study) and when sample sizes were reduced to 100 the rule usually gave correct results. Apparently 4 different simple structure types

14

were examined over some 160 randomly generated sample
correlation matrices. Hamburger fails to report the
number of population common factors and variables that
were used to generate his sample correlation matrices
and as he points out, "These conclusions are of course
tempered by the fact that results probably depend on
the number and size of common factors and the number of
variables ....".

Joreskog (1963) studied the sampling errors of
individual loadings on unrotated common factors. He
used a factor analytic technique which he developed in
the 1962-1963 time period (Joreskog,1962,1963).  this
procedure is reported to yield factors very similiar in
appearance to those generated by the PCA procedure.  In
this particular study Joreskog generated a small number
of sample correlation matrices, under various
conditions, and analyzed various characteristics of
them.  Two cases stand out in particular. In one case,
six 20-variable sample correlation matrices were
generated over uncorrelated factors. In three of the
matrices, the factor scores were assumed to be normally
distributed while in the other three matrices the
factor scores were assumed to have a rather strong
skewness.  Sample sizes for each set of three
correlation matrices were 100, 200, and 300. When the
sample correlation matrices were factor analyzed to
yield sample factor loading matrices it was found that

15

the root mean square (rms) deviation of sample loadings to population loadings was somewhat less than $1/\sqrt{N}$, where N is sample size, the approximate standard error of a zero correlation. There were no consistent differences reported due to skewness.

A second case showed a sharp contrast to the first case. The population factor structure chosen for examination, in this case, was a 10-variable three factor structure. The structure exhibited perfect simple structure. Joreskog generated 10 sample correlation matrices for each of the three sample sizes. This time the differences between the population loadings and the sample loadings were very great. In some samples one or more of the factors generated could not be confidently matched with population factors. These errors decreased slightly as the sample size was increased.

Joreskog next decided to rotate the previous set of factors, via a least squares procedure, to its population structure. The resultant standard errors of the rotated loadings were quite small; somewhat smaller than $1/\sqrt{N}$. Also, the sampling errors of non-zero loadings tended to be smaller than those for zero loadings, in the same manner that sampling errors for the Pearson correlation coefficients are proportional to $1-r^2$. Browne (1968) also observed that the sampling errors of rotated factor loadings were about the same

as correlation coefficients, of the order $1/\sqrt{N}$.

Joreskog also looked at the sampling error of rotated loadings when the factor scores were given from a rectangular distribution. He found that although the standard errors of the loadings rose slightly in this case they remained less than $1/\sqrt{N}$ and exhibited the same proportionality to the magnitude of the original loadings as found in the normally distributed factor score cases. This result (coupled with the skewed factor score distribution results previously mentioned) suggest that factor analytic methods are reasonably robust in respect to moderate departures from the normally distributed factor score assumption.

Hamburger (Cliff and Hamburger, 1967) observed a bias in the estimates of individual factor loadings. Hamburger does not state his original factor structure but he does state that the sample sizes studied were 100 and 400. He noticed that larger loadings (.6 to .9) were consistently underestimated while the smaller loadings (.2 to .5) showed a tendency to be overestimated. Browne's (1968) data exhibited a similiar trend. The factor analytic procedure used for these cases was PCA with squared multiple correlations placed on the diagonal of the input correlation matrix.

Cliff and Pennell (1967) studied in some detail the effects of communality, factor strength, and loading size on the sampling characteristics of factor

loadings. The sampling characteristics addressed in this study were referred to as "stability" and "bias". Stability refers to the amount of variability an individual factor loading might be expected to exhibit from sample to sample. Bias refers to the extent to which the mean of a sampling distribution of factor loadings might be expected to approximate the population loading.

Two model population factor structures were studied. Each factor structure was constructed with four different factor strengths, four different communalities, and four different loading sizes. The loadings were situated in such a fashion to facilitate various comparisons. For instance, by choosing selected loadings, it was possible to compare the effect of factor strength on given loadings of equal magnitude and communality. The loadings of the first model structure ranged from .9 to .45, while the second's loadings ranged from .7 to .35. Fifty sample correlation matrices were generated for each of the model structures. The sample sizes were not given. These matrices were factor analyzed by the PCA procedure with communalities placed on the diagonal. Four factors were extracted and the resultant structure was rotated, via a least squares procedure due to Cliff (1966), to fit the model structure. The means and standard deviations of the individual factors were

18

calculated over the 50 sample structures. Individual
factor loadings exhibited a wide variety of sample
frequency distributions.

Several interesting influences affecting stability
were noted. A non-zero loading which was associated
with a large communality was almost always associated
with smaller standard deviations then those loadings
associated with smaller communalities. A similiar trend
was noticed for zero loadings. It was also noticed that
the larger a loading was, the smaller its standard
deviation tended to be. The trends were presented quite
clearly in graphs in the text of the article. The
authors further noticed that stronger factors tend to
produce more stable loadings, apparently independent of
the other parameters.

Cliff and Pennel next summarized these
observations in a multiple regression study. The
dependent variable was the standard deviation of the
factor loadings. There were 7 independent variables::

   1. Size of population factor loading
   2. Population communality of the variable
   3. Number of non-zero loadings on the factor
   4. Number of non-zero loading on the variable
   5. Squared loading on the factor
   6. Discrepancy between loading and other loadings
on the factor
   7. Mean loading on the factor


The correlation matrix was prepared and
correlations of the order .80 were noted for predictors
1,2,5,7. Using a stepwise regression routine it was

19

found that six of the variables could account for
nearly 89% of the variance present. The best two
predictors were 6 and 7 which together explained 86% of
the variance. The authors, however, point that
predictor pairs 2 and 5 or 2 and 7 give $r^2$s of .842 and
.843 respectively. No power transformations or
interactions were tested. The authors also report
instances of bias in the sample factor loading matrix
but do not attempt to characterize it.

Cliff and Pennell summarize by concluding

"...that communality rather than loading size
is the important determiner of
stability...Higher communalities mean not
only greater stability for the loadings of
specific tests (variables) but also load to
stronger factors which mean that the
stability of all the loadings is improved."

Cliff and Pennell did not address sample size and
its possible interaction with other effects. Also the
number of variables and factors were not varied.

Pennell (1968) extended the study by looking at
the influences of communality and N (N is the sample
size used to generate sample correlation coefficients)
on the sampling distributions of factor loadings.
Pennell inserted variables with different communalities
in randomly constructed factor structures. Samples were
drawn from the population correlation matrix to prepare
a sample correlation matrix. The sample correlation
matrix was then factor analyzed using PCA with squared
multiple correlations on the diagonals. As in Cliff and

20

Pennell's 1967 study, the standard deviation of the factor loading was taken to be the dependent variable.

Pennell noticed in the 1967 study that univocal variables (variables which load on a single factor) not only facilitated rotation and subsequent factor interpretation but also resulted in smaller sampling errors. Further it was noticed that zero loadings seem to exhibit a greater degree of variability than a non-zero loadings. Pennell hypothesised that the variability of a loading might increase with its complexity across the factors. To avoid the confounding of error in complex variables Pennell used only univocal tests inserted in randomly constructed factor structures.

The research design chosen was a two way analysis of variance with 5 levels of the two independent variables, N and communality. The levels of N were taken as 100, 150, 300, 600, and 2500. The levels of communality were taken as .1, .3, .5, .7 and .9. for the test variables. The test variables were inserted into a randomly constructed, 12 variable by 2 factor, population structure (The test variable's position in the structure was also determined in a random fashion.) Three such random structures were generated for each of the ANOVA's 25 cells. These 75 structures then were used to generate 100 sample correlation matrices. These correlation matrices were then factor analyzed in the

21

fashion previously mentioned and rotated via Cliff's
(1966) procedure back to the population factor
structure. The standard deviations of the zero and
non-zero sample correlation matrices were calculated
for both blocks of correlation matrices. Hence each
cell of the ANOVA contained three replications of
sample standard deviations. The subsequent fixed
effects, two way ANOVA revealed that both main effects
and their interactions were significant. This was true
for both non-zero and zero loadings. It was noticed
that while the F ratio for N remained approximately the
same for both non-zero and zero loadings, the F ratio
due to communality was strongest for non-zero loadings.

Graphs were drawn that depicted the standard
deviations (non-zero loadings) of the various test
variables as functions of $1/\sqrt{N}$. The results were
striking, in all but one case, clear linear trends were
observed. It was also clear from the graphs that the
standard deviations were conditional on the magnitude
of the variable's communality. Similar trends were
observed for zero loadings although differences due to
communalities were harder to discern.

Pennell asserts that his study has demonstrated
clearly the advantage of developing factor pure
(univocal) variables for studies of psychological
traits, because these variables exhibit the smallest
amount of sampling errors in the non-zero loadings (the

22

loading usually of most interest.)

Pennel presents tables which define 95% confidence
intervals about the zero loadings as a function of N.
These tables graphically refute the rule of thumb that
selects only loadings greater than .3 as being
significant. This author reminds the reader that a
"forced fit" rotation scheme was employed to generate
these values and as such the values in the table arise
when the correct structure is known, a priori, and the
sample loading structure is rotated to fit it. Hence
errors here are due to the factor analytic method and
sampling errors. Obviously, investigators may imploy
factor analytic and rotation schemes which will produce
effects other than those due only to sampling error.
Pennell closes his paper with an interesting table in
which he shows the size of sample loadings necessary to
be significately different from the non-zero loadings
he tested (.9, .8, .7, .6, .5) at an alpha of .05. For
example, when N is 100, a sample loading would have to
be less than .79 to reject the hypothesis that it was
actually .90.

Manners and Brush (1979) studied the "reliability"
of factor analytic techniques. Reliability is defined
as (a) the mean squared error between factor loadings
for sample and population factor loading structures and
(b) the ability of a factor analytic model to capture
correctly the number of factors in the population. The

23

research endeavored to compare the reliability of four
separate factor analytic techniques with respect to the
effects of sample size, number of variables, and number
of factors. An analysis of variance approach was used
where the treatments were taken as a) the number of
variables, b) the number of observations (observations
on the variables), and c) the four different factor
analytic models. All possible interactions were also
examined.

The experimental procedure called for use of a
factor structure due to Browne (1968). This structure
was divided into two experimental conditions; the first
was 16 variables and 4 factors and the second was 12
variables with the same 4 factors. (The reader should
note that factors II and III in Browne's structure are
not orthogonal. The angle between these factors is
aproximately 77 degrees. Since only an orthogonal
rotation was used in the analysis, one wonders why
Browne did not employ only mutually orthogonal
factors.) Ten random correlation matrices were
generated from each of the experimental population
factor loading structures for sample observation sizes
of N=100 and N=500, respectively. Hence 10 times 2
times 2 = 40 sample correlation matrices were factor
analyzed by the four differnent techniques. The four
different techniques included PCA with initial
communality estimates placed on the diagonal.

24

Three dimensionality assessment rules were tested for the hybrid PCA technique. The first two rules were to choose the number of factors associated with the eigenvalues of magnitude greater than 1 and 0, respectively. Saunder's method which was mentioned earlier in the literature review, was also tested. For the eigenvalues greater that 1 rule, 28% of the dimensionality assessments were correct. The remaining assessments were within 2 factors of being correct with 35% predicting 3 factors and 37% predicting 2 factors. The eigenvalues greater than 0 rule was within one factor for 45% of the assessments and high for the rest, and the variance was much larger for this rule than the eigenvalues greater than 1 rule.

In the fixed effects ANOVA experiment each cell contained 10 replications of mean square loading errors (10 correlation matrices for each cell; 4 (models) times 2 (sample sizes) times 2 (# of variables) = 16 cells). All the treatments and treatment interactions proved significant (alpha=.05), save the interaction between number of variables and number of observations.

To summarize; Manners and Bush provide evidence that factor analytic reliability is influenced by

1) The specific factor analytic models chosen.
2) The interaction of factor analytic model choice and number of variables.
3) The interaction of factor analytic model choice and sample size.
4) The interaction of factor analytic model choice, sample size and number of variables.

25

In passing, one also notes that all the main treatments: variables, observations, and models were significant. Sampling error decreased as observations increased. Sampling error also decreased as the number of variables increased, much as is observed in multiple regresion analysis.

This ends the literature review section dealing with Monte Carlo experimental work in factor analysis. The following paragraphs present a brief overview of recent literature which is related either to the research in this report or the factor analysis problem in general.

The rotational scheme used in this report is due to Schoneman (1966). This rotational procedure allows one to rotate a sample factor loadings matrix to given target matrix (usually the hypothesized population factor loadings matrix). The rotation is accomplished in such a manner as to minimize, in a least squares sense, the residual differences between the rotated matrix and its target.

Odell and Feiveson (1966) provide the methodology by which all the reviewed studies generated sample values from multinormal populations with given covariance structures. An algorithm for the bivariate normal is given by Naylor (1966).

A mathematical entity known as the condition number of a matrix is used in this report. The

condition number of a matrix is an especially useful tool in systems of linear equations. If small errors in the right-hand side or coefficients of a linear system produce a large effect on the solution, then the system of equations is said to be ill-conditioned. The condition number of a matrix serves as an index of this ill-conditioning. The condition number of a matrix is given by the absolute value of the ratio of that matrix's largest eigenvalue to its smallest. Westlake (1968) offers a clearly written text on the application of the condition number and other measures as applied to matrix inversion and linear equations. Belsley, et al. (1980) discuss the effects of ill conditioning by applying the condition number in detecting collinearity in multiple regression.

Anderson (1958) is recommended for rigorous yet succinct theoretical treatment of PCA. Joreskog (1979) offers a novel introduction to factor analysis and its associated vernacular using the concept of partial correlations as a starting point. The second paper in the book deals with statistical tests for confirmatory factor analysis, the only such statistical tests this author was able to find. Also there is an interesting article by Catell and Sullivan (1962) in which the concepts of factor analysis are made clearer to the novice through the use of a physical example using cups of coffee.

27

## Objectives of the Research

The objectives of this research effort were the following:

1) Develop software which will allow a user to study the influences of sample size, number of variables, number of factors, and complexity of the factor structure in PCA. This software was to allow a user to input either a particular structure or given covariance matrix.

2) Address the three questions raised in the problem statement.

3) Summarize recent developments in this area as found during a literature search.

## Scope of the Research

**Dimensionality Assessment.** Kaiser's criterion and
the Scree test are the two dimensionality assessment
procedures to be addressed in this report.  These two
procedures appear to be the most popular.  It was also
felt that testing Horn's procedure on the CDC 6000
would prove cost prohibitive.

**Factor Interpretation.**  This report did not treat
the problem of selecting the most appropriate
rotational technique.  This report assumes that the
correct rotational scheme is applied.  This report
attempts to provide a tool by which investigators can
estimate the magnitude of factor loading errors to be
incurred under various experimental conditions.  In
particular, sample sizes were taken from the range
10-100, the number of variables ranged 10-15, the
number of factors ranged from 3-6. In total 27 separate
population structures were examined. Each structure was
examined at sample sizes of 10, 25, 50, 100.

## Approach to the Research

The research design used in this report was as follows: Initially, seven 10 variables by 3 factors structures were examined. These structures will be referred to as the "original" structures. The original structures are given in figure 1. These original structures are the lower left hand point in the graph of the research design given in figure 2, with coordinate (3,10). Next, four of the original structures were selected to be perturbed by the presence of added nuisance variables. These structures are given in figure 3. Five and two nuisance variables were added to these structures,respectively. The nuisance variables were chosen to load on single factors and have communalities between .01 and .09. These eight structures are coordinates (3,15) and (3,12) on the research design graph. The coordinates (4,10) and (6,10) consist of three structures each, perturbed by the presence of one and three nuisance factors. The nuisance factors were constructed such that all factors were orthogonal (save factors IV and VI in structure 17, they form an angle of approximately 74 degrees, as it turned out this discrepancy was insignificant, as in Browne's (1968) structure). These structures are given in figure 4. To complete the design the same three structures were perturbed by both nuisance variables and factors. These structures are

30

```
       1              2              3              4

1 0 0         .9  0  0      .8  0  0      .7  0  0
0 1 0          0 .9  0       0 .8  0       0 .7  0
0 0 1          0  0 .9       0  0 .8       0  0 .7
1 0 0         .9  0  0      .8  0  0      .7  0  0
0 1 0          0 .9  0       0 .8  0       0 .7  0
0 0 1          0  0 .9       0  0 .8       0  0 .7
1 0 0         .9  0  0      .8  0  0      .7  0  0
0 1 0          0 .9  0       0 .8  0       0 .7  0
0 0 1          0  0 .9       0  0 .8       0  0 .7
1 0 0         .9  0  0      .8  0  0      .7  0  0
```

```
       5              6              7

.9  0  0      .7 .7  0      .7 .7  0
 0 .8  0      =7 .7  0      =7 .7  0
 0  0 .7       0  0 .9       0  0 .6
.6  0  0      .8  0  0      .5  0  0
 0 .5  0       0 .8  0       0 .4  0
 0  0 .4       0  0 .7       0  0 .3
.3  0  0      .6  0  0      .6  0  0
 0 .2  0       0 .7 =7       0 .7 .7
 0  0 .7       0 .7 .7       0 .7 =7
.7  0  0      .7  0  0      .7  0  0
```

Figure 1.    Original Factor Structures

Figure 2. Research Design

```
       8              9              10             11

   1  0  0      .9  0  0      .7 .7  0      .8  0  0
   0  1  0       0 .9  0      ÷7 .7  0       0 .8  0
   0  0  1       0  0 .9       0  0 .9       0  0 .8
   1  0  0      .9  0  0      .8  0  0      .8  0  0
   0  1  0       0 .9  0       0 .8  0       0 .8  0
   0  0  1       0  0 .9       0  0 .7       0  0 .8
   1  0  0      .9  0  0      .6  0  0      .8  0  0
   0  1  0       0 .9  0       0 .7 ÷7       0 .8  0
   0  0  1       0  0 .9       0 .7 .7       0  0 .8
   1  0  0      .9  0  0      .7  0  0      .8  0  0
   0 .2  0       0 .2  0       0 .2  0       0 .2  0
   0  0 .1       0  0 .1       0  0 .1       0  0 .1
  .3  0  0      .3  0  0      .3  0  0      .3  0  0
   0 .2  0       0 .2  0       0 .2  0       0 .2  0
   0  0 .1       0  0 .1       0  0 .1       0  0 .1


      12             13             14             15

   1  0  0      .9  0  0      .7 .7  0      .8  0  0
   0  1  0       0 .9  0      ÷7 .7  0       0 .8  0
   0  0  1       0  0 .9       0  0 .9       0  0 .8
   1  0  0      .9  0  0      .8  0  0      .8  0  0
   0  1  0       0 .9  0       0 .8  0       0 .8  0
   0  0  1       0  0 .9       0  0 .7       0  0 .8
   1  0  0      .9  0  0      .6  0  0      .8  0  0
   0  1  0       0 .9  0       0 .7 ÷7       0 .8  0
   0  0  1       0  0 .9       0 .7 .7       0  0 .8
   1  0  0      .9  0  0      .7  0  0      .8  0  0
   0 .2  0       0 .2  0       0 .2  0       0 .2  0
   0  0 .2       0  0 .2       0  0 .2       0  0 .2
```
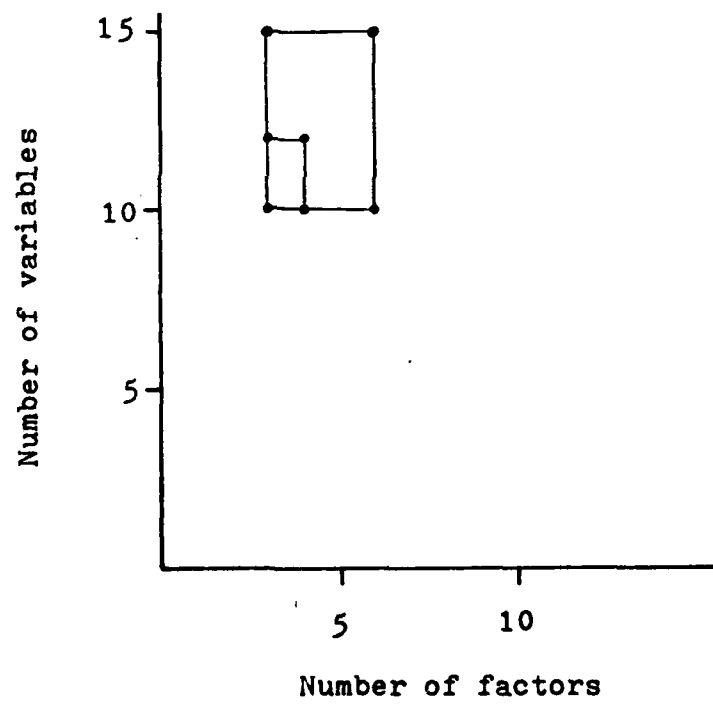
Figure 3.   Perturbed Factor Structures: Nuisance Variables

33

**16**

```
.9 0 0 .2 0 0
0 .9 0 0 .2 0
0 0 .9 0 0 .2
.9 0 0 -2 0 0
0 .9 0 0 -2 0
0 0 .9 0 0 -2
.9 0 0 0 0 0
0 .9 0 0 0 0
0 0 .9 0 0 0
.9 0 0 0 0 0
```

**17**

```
.7 .7 0 0 0 0
.7 .7 0 0 0 0
0 0 .9 0 .2 0
.8 0 0 -15 0 -14
0 .8 0 0 0 0
0 0 .7 0 -.257 0
.6 0 0 .2 0 -0105
0 .7 .7 0 0 0
0 .7 .7 0 0 0
.7 0 0 0 0 .25
```

**18**

```
.8 0 0 .2 0 0
0 .8 0 0 .2 0
0 0 .8 0 0 .2
.8 0 0 0 0 0
0 .8 0 -2 0 0
0 0 .8 0 -2 0
.8 0 0 0 0 -2
0 .8 0 0 0 0
0 0 .8 0 0 0
.8 0 0 0 0 0
```

**19**

```
.9 0 0 .2
0 .9 0 0
0 0 .9 0
.9 0 0 -2
0 .9 0 0
0 0 .9 0
.9 0 0 0
0 .9 0 0
0 0 .9 0
.9 0 0 0
```

**20**

```
.7 .7 0 0
-7 .7 0 0
0 0 .9 .2
.8 0 0 0
0 .8 0 0
0 0 .7 -257
.6 0 0 0
0 .7 -7 0
0 .7 .7 0
.7 0 0 0
```

**21**

```
.8 0 0 .2
0 .8 0 0
0 0 .8 0
.8 0 0 -2
0 .8 0 0
0 0 .8 0
.8 0 0 0
0 .8 0 0
0 0 .8 0
.8 0 0 0
```

Figure 4.    Perturbed Factor Structures:
Nuisance Factors

34

```
        22                        23                          24

.9 0 0 .2 0 0          .7 .7 0 0  0   0        .8 0 0 .2 0 0
0 .9 0 0 .2 0          .7 .7 0 0  0   0        0 .8 0 0 .2 0
0 0 .9 0 0 .2          0 0 .9 0 .2   0         0 0 .8 0 0 .2
.9 0 0 .2 0 0          .8 0 0 .15 0  .14       .8 0 0 0 0 0
0 .9 0 0 .2 0          0 .8 0 0  0   0         0 .8 0 .2 0 0
0 0 .9 0 0 .2          0 0 .7 0 .257 0         0 0 .8 0 .2 0
.9 0 0 0 0 0           .6 0 0 .2 0  .0105      .8 0 0 0 0 .2
0 .9 0 0 0 0           0 .7 .7 0  0   0        0 .8 0 0 0 0
0 0 .9 0 0 0           0 .7 .7 0  0   0        0 0 .8 0 0 0
.9 0 0 0 0 0           .7 0 0 0  0   .25       .8 0 0 0 0 0
0 .2 0 0 .3 0          0 .2 0 0  0   0         0 .2 0 0 .3 0
0 0 .1 0 0 0           0 0 .1 0  0   0         0 0 .1 0 0 0
.3 0 0 0 0 0           .3 0 0 0  0   0         .3 0 0 0 0 0
0 .2 0 0 .3 0          0 .2 0 0 .3   0         0 .2 0 0 .3 0
0 0 .1 0 0 0           0 0 .1 0  0   0         0 0 .1 0 0 0


        25                    26                      27

.9 0 0 .2              .7 .7 0 0              .8 0 0 .2
0 .9 0 0              .7 .7 0 0              0 .8 0 0
0 0 .9 0              0 0 .9 0              0 0 .8 0
.9 0 0 .2              .8 0 0 0              .8 0 0 .2
0 .9 0 0              0 .8 0 0              0 .8 0 0
0 0 .9 .05            0 0 .7 .0643          0 0 .8 .0563
.9 0 0 0              .6 0 0 0              .8 0 0 0
0 .9 0 0              0 .7 .7 0              0 .8 0 0
0 0 .9 0              0 .7 .7 0              0 0 .8 0
.9 0 0 0              .7 0 0 0              .8 0 0 0
0 .2 0 0              0 .2 0 0              0 .2 0 0
0 0 .2 .225          0 0 .2 .225          0 0 .2 .255
```

Figure 5.    Perturbed Factor Structures:
            Nuisance Factors and Variables

given in figure 5. These structures are coordinates (4,12) and (6,15) in the research design graph.

Each of the 27 structures were examined over sample sizes of 10,25,50, and 100. These sample sizes are considerablely lower than those used in the reviewed studies: 100 and 400 were used by Hamburger (1967), Joreskog (1963) used 100,200, and 300, Cliff and Pennell (1968) used 100,150,300,600, and 2500, while Manners and Brush (1979) used 100 and 500. These studies reported that experimental results were, in general, only moderately improved by increasing sample size. Therefore, it was thought that the lower sample sizes would provide not only more interesting results but shed some light on just how many samples are required to perform an accurate PCA in the given experimental region.

Each structure and sample size combination was analyzed according the following experimental procedure:

1) A population covariance matrix was formed by multiplying the population structure matrix by its transpose.

2) The appropriate number of sample vectors were drawn randomly from the population covariance matrix.

3) The sample vectors were then used to form a sample correlation matrix. The condition number of this matrix was calculated at this step.

36

4)   The sample correlation matrix was then factor analyzed by the PCA procedure.  The correct number of factors were retained.  Dimensionality assessment statistics were collected at this step.

5)   A factor loadings matrix was prepared and rotated via a least squares procedure due to Schoneman (1966) back to the original population structure.  The mean square deviation of sample loadings from population loadings was calculated at this step.

6)   Steps 2-5 were repeated 1000 times for each structure-sample size combination.

## Complications in a Population Structure

Rationale. One difficulty in a study such as this is in determining what structures should be examined. The limited experience of this author indicates that "simple" structures, whose variables load on no more than a few factors and which contain many zeros, are probably of the greatest use. This intuitive feel is merely a vague generalization of Thurstones criteria (see Harman, 1967, pg. 98) which is widely accepted as a desirable quality of a population structure. In any event, more "complicated" structures may be very difficult, if not impossible, to give any meaningful interpretation to.

Once one has decided to use these simple structures one might wish to study several such structures. The question arises, how does one compare different factor structures? It stands to reason that a perfect, simple population structure (all ones for loadings with each variable loading on a single factor) will be easier to detect from sample data than a structure with low factor loadings and variables which load on more than one variable. Here the second structure could be said to be more complicated than the first. It would be desirable to have an index number which could be derived from a given structure. This index number should grow in magnitude as a structure becomes increasingly complicated. One such candidate

38

is the average uniqueness of the structure. Pennell
(1968) rejects such a measure because it does not take
into account the fact that variables of equal
communality may load on differing numbers of factors.

A given factor structure is said to be more
complicated than another if the first factor structure
is harder to glean from experimental data than the
second. The following is a proposed index for the
complexity of a population structure.

Complexity Index. It was reasoned that
complications in a given structure are due to two
components:

1) Complication due to structure - if
manisfestation variables load on single variables then
a simple structure exists. As the manisfestation
variables load significantly on more than one factor
the complexity of the structure increases.

2) Complication due to Uniqueness - if
manisfestation variables demonstrate high communalities
across the factors (low uniqueness) then there is a
higher chance of closely reproducing this factor
structure than that of another structure with higher
uniqueness.

Let the complexity index be defined by the
quantity:

$$\frac{\sum\limits_{i=1}^{M} \sum\limits_{j=2}^{N} \sum\limits_{k=1}^{J-1} (a_{ik}a_{ij})^2}{M} + \left( 1 - \frac{\sum\limits_{i=1}^{M} h_i}{M} \right)$$

where the $A_{ij}$ are factor loadings in the $ij$th position.
N is the number of factors, and M is the number of
variables.  The $H_i$ are the communalities of the ith row
of the factor structure matrix.

The first term is the complexity due to structure.
The second term is the complexity of the structure due
to uniqueness.

The second term is simply the average uniqueness of
the structure.  As the average uniqueness grows the
complexity grows.  This term is bounded by 0 and 1.

The first term is the quartimax criterion divided
by the number of variables, M.  The quartimax criterion
is minimized at 0 when perfect simple structure is
present.  As variables start to load on more than one
factor this quantity grows.  In order for this quantity
to be useful as a component of an index it has to be
bounded.  The lower bound is zero. The upper bound can
be found if the following maximization problem can be
solved.

40

MAXIMIZE

$$\sum_{i=1}^{M} \sum_{j=2}^{N} \sum_{k=1}^{J-1} (a_{ik} a_{ij})^2$$

SUCH THAT

1) $|a_{ij}| \leq 1$ $\forall i = 1,\ldots,M,$ $\forall j = 1,\ldots,N$

2) $\sum_{k=1}^{N} (a_{ik})^2 \leq 1$ $\forall i = 1,\ldots,M$

3) $\sum_{i=1}^{M} a_{ik} a_{ij} = 0$ $\forall k \neq j$

The first constraint merely requires the loadings to be less than or equal to unity. The second constraint requires each variable's communality to be less than or equal to unity. The third constraint reflects the mutual orthogonality of the factors.

If only the first constraint is taken to be binding an upper bound of

$$M \cdot \binom{N}{2}$$

can be established by setting all the elements in the matrix to 1.

If we require both the first and second constraint to be binding then:

$$\sum_{i=1}^{M} \sum_{j=2}^{N} \sum_{k=1}^{J-1} (a_{ik} a_{ij})^2 \leq \sum_{i=1}^{M} \sum_{j=2}^{N} \sum_{k=1}^{N} (a_{ik})^2 (a_{ij})^2$$

41

and

$$\sum_{i=1}^{M} \sum_{j=2}^{N} \sum_{k=1}^{N} (a_{ik})^2 (a_{ij})^2 = \sum_{i=1}^{M} \sum_{j=2}^{N} (a_{ij})^2 \sum_{k=1}^{N} (a_{ik})^2$$

and

$$\sum_{i=1}^{M} \sum_{j=2}^{N} (a_{ij})^2 \sum_{k=1}^{N} (a_{ik})^2 \leq \sum_{i=1}^{M} \sum_{j=2}^{N} (a_{ij})^2 \cdot 1$$

similarly:

$$\sum_{i=1}^{M} \sum_{j=2}^{N} (a_{ij})^2 \cdot 1 \leq \sum_{i=1}^{M} 1 \cdot 1 = M$$

hence M also acts as a weak upper bound.


This author was not able to determine a upper bound with the third constraint binding.

To summarize, the above index is submitted as a possible candidate to compare differing structures for inherent complexity. This index does not attempt to account for the possible influence due to the ratio of the number of variables to the number of factors. In light of this fact, comparisons in this report using the complexity index are only done across structures with the variable to factor ratio held constant. Two other points should be noted. First, a weak upper bound is used to normalize the structural complexity term. Undoubtedly a stronger upper bound computed with the third constraint binding would lead to a stronger

42

index.  Secondly, it was assumed that the weights on
the two terms were equal.  This assumption implies that
complexity is equally attributed to structure and
uniqueness.  If the index proves promising, perhaps
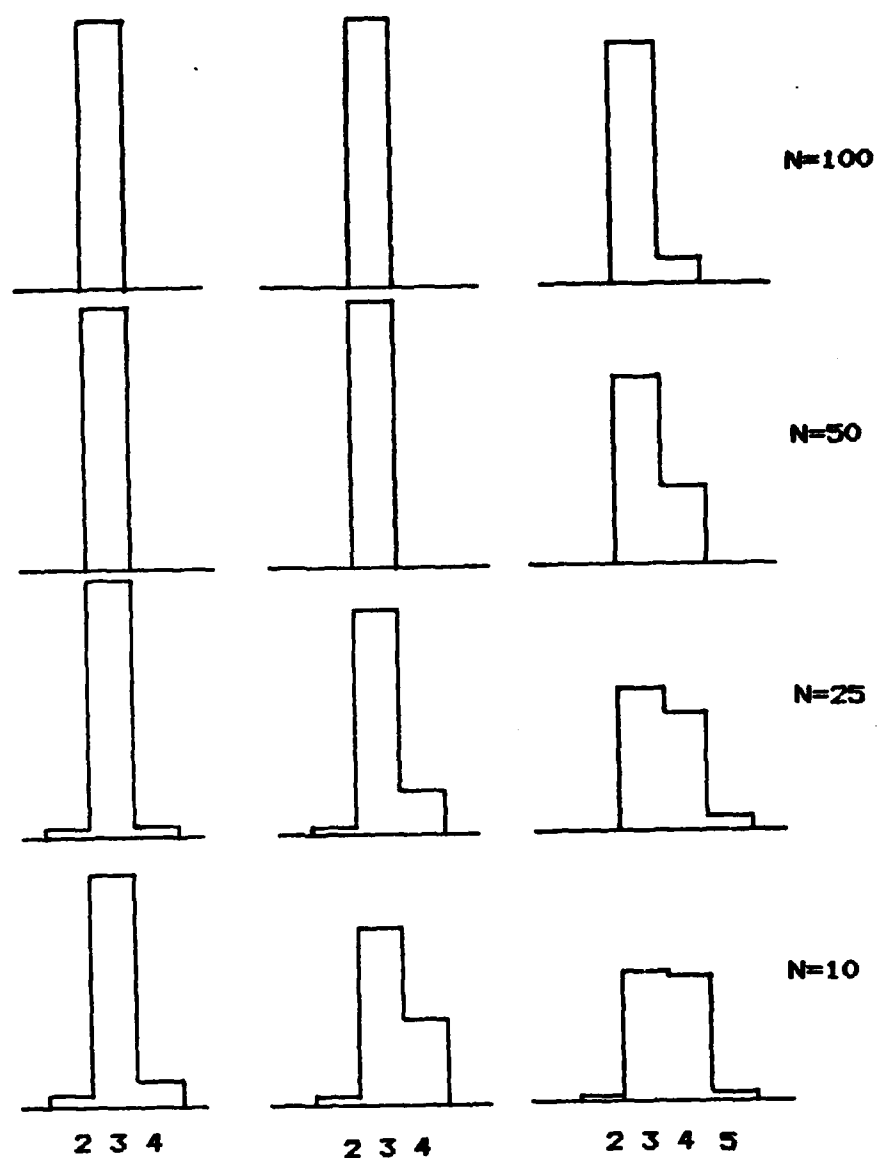future regression studies could address this issue.

Dimensionality Assessment Analysis

The two dimensionality assessment techniques addressed in this section are Kaiser's criterion and Catell's scree test.

Kaiser's criterion is to merely retain all factors whose associated eigenvalues ae greater than or equal to 1. Catell's scree test is a graphical technique in which an investigator looks for a break in a plot of rank ordered eigenvalues. This section does not attempt to make statistical statements about these techniques.

Kaiser's Criterion. For all the structures examined in this study, all dimensionality assessments based on Kaiser's criterion were within two factors of being correct. This is probably attributed to the structurally "clean" sampling populations studied and the low factor to variable ratios used. Most dimensionality assessments were, in fact, within one factor of being correct.

Histograms of dimensionality assessments due to Kaiser's criterion are given in figure 6 for structures 1,3, and 7. The correct dimensionality for each structure is 3. Each histogram contains a total of 1000 dimensionality assessments. Note that the variability in the dimensionality assessments is markedly larger in structure 7 than in structure 3. Although the percentage difference in average

44

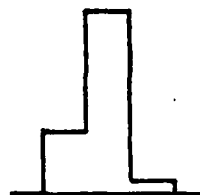| Structure: | 1 | 3 | 7 |
|---|---|---|---|
| Avg. Comm.: | 1.0 | .64 | .563 |
| Complexity: | 0 | .36 | .533 |

Figure 6. Dimensionality Assessment Histograms

using Kaiser's Criterion

communality is only about 14%, the percentage difference in the complexity index is 48%. No solid conclusions can be drawn at this point, but at least one notices that the complexity index is moving in the right direction.
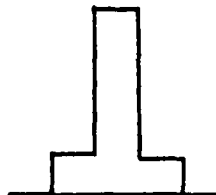
In figure 7 histograms are presented for structure 23. This structure is perturbed by 3 nuisance variables and 3 nuisance factors. It is one of the more complicated structures stuidied in this report. Although, in an absolute sense there are 6 inherent population factors to this structure, 3 of these factors account for less than 3% of the total variance which could be explained by this structure. As can be seen in the histograms a dimensionality assessment of 3 was never made. Clearly the addition of nuisance factors and variables can impact dimensionality assessments via Kaiser's criterion. However, for the structures studied here, one can expect to be within 2 factors of the true dimensionality. To this author, Kaiser's criterion seems to be a good rule for thumb for dimensionality assessment.

Catell's Scree Test. Since this test is graphical in nature, it was very difficult to conceive of a method to apply Mr te Carlo techniques to its analysis. Clearly, one could not hope to examine a thousand graphs visually within a limited time period.
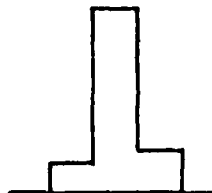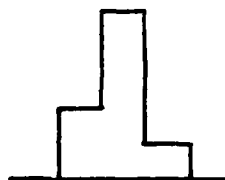
Catell's scree test is a graphical technique used

N=100

N=50

N=25

N=10

4 5 6

Structure:    23

Avg. Comm.:  .532

Complexity:  .542

Figure 7. Dimensionality Assessment Histograms

to visually locate the hypothesized break in ranked

eigenvalue magnitudes which should occur just before

that eigenvalue which is associated with the correct

dimensionality. The Scree test is explained in the

literature review section of this thesis.

Catell would have an investigator retain factors

down to and including the factor which begins his scree

line. To test this procedure the following approach

was taken. If the scree test is an acceptable

procedure then certainly one would expect the method to

work well under ideal conditions. An ideal condition

for an investigator would occur if he were sampling

from a population like structure 1. Figures 8,9,10,

and 11 are Catell's scree test for sample sizes

10,25,50, and 100 respectively for structure 1. Each

plotted point is the mean of the ith ranked eigenvalue

over 1000 trials at the particular sample size.

Approximate 95% confidence intervals are provided for

the means of the eigenvalues at the correct

dimensionality (3, in this case) and 1 plus the correct

dimensionality. Note in figure 8 that there is no

apparent break in the means of the ranked eigenvalues.

This situation improves markedly as the sample size

increases, figures 9,10, and 11. Notice in figure 11

that a definite break in magnitude is present between

eigenvalues 3 and 4. Further, notice that the

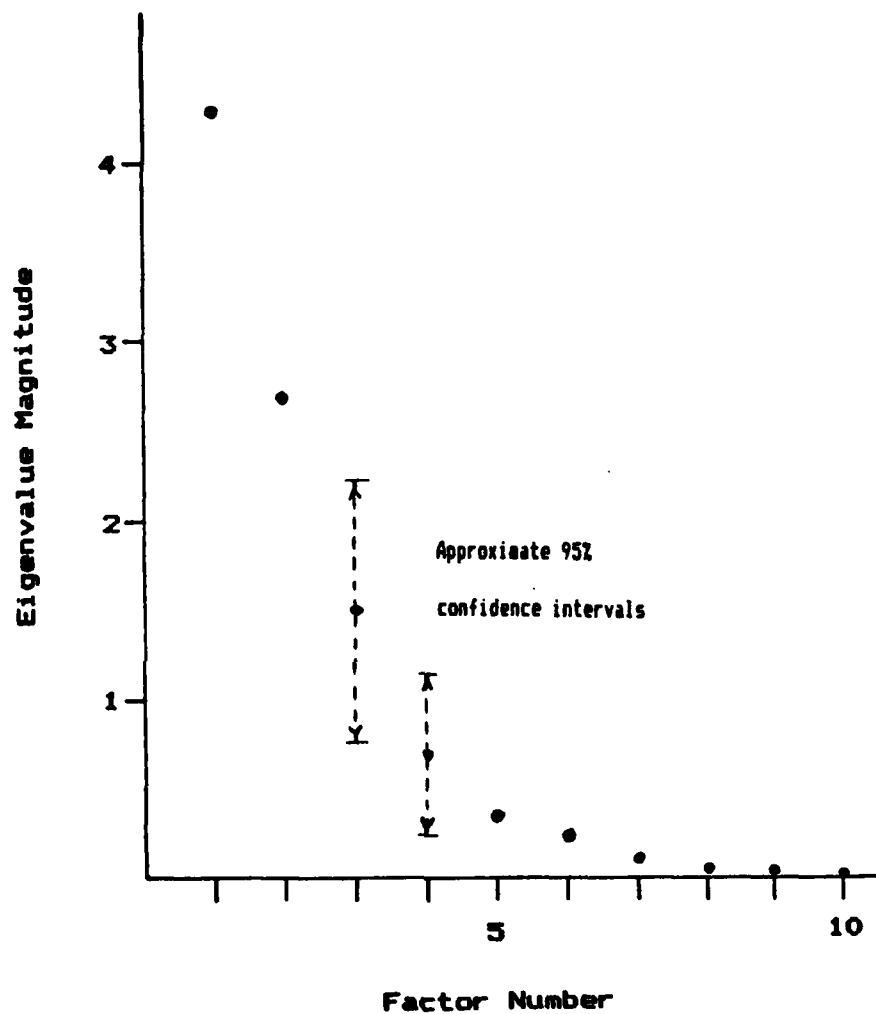confidence intervals for the two eigenvalues do not
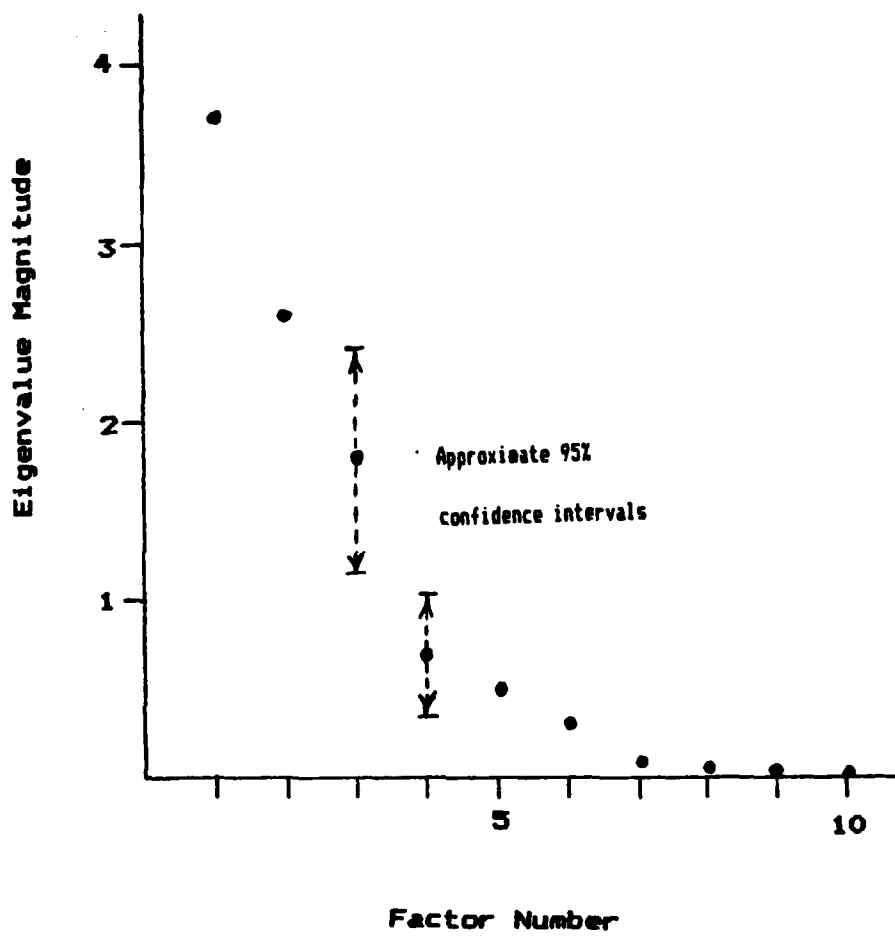
Figure 8. Scree Test, Structure 1, N=10

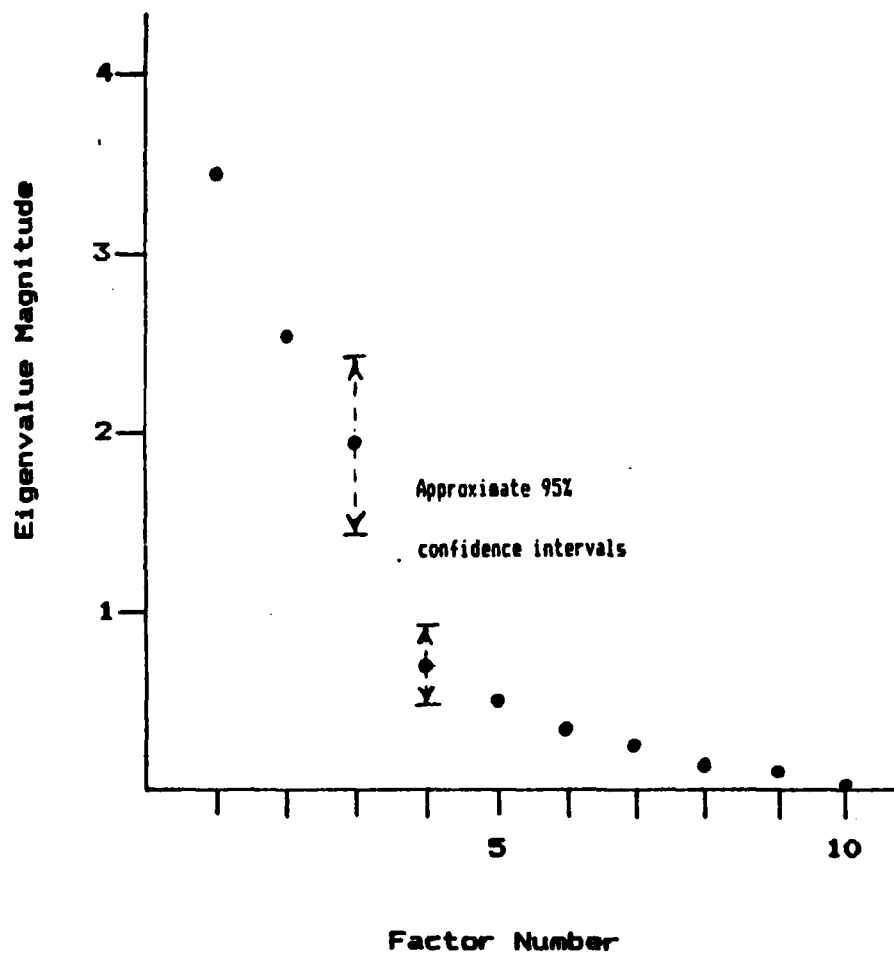Figure 9. Scree Test, Structure 1, N=25
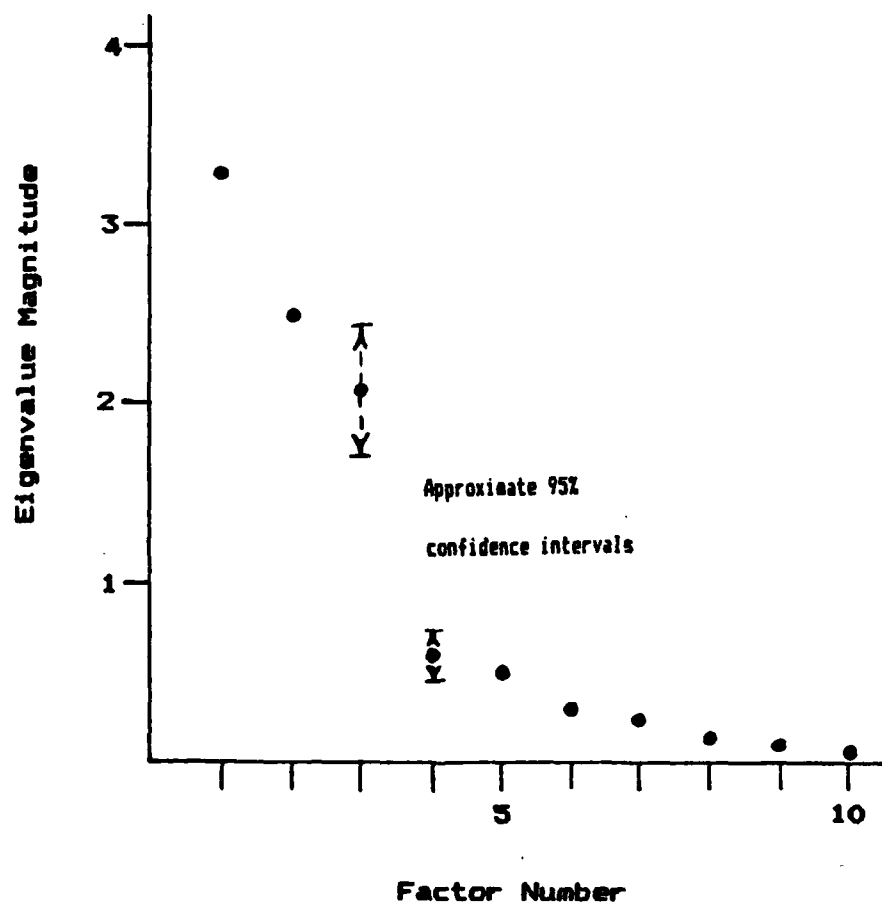
50

Figure 10. Scree Test, Structure 1, N=50

51

Figure 11. Scree Test, Structure 1, N=100

52

overlap. If one were to apply Catell's scree test to
the means of these eigenvalues, clearly, one would
retain 4 factors. Thus, when sampling is accomplished
under even ideal conditions Catell's test has yielded
incorrect results. In fairness to Catell, however,
figure 11 could be said to exhibit what Catell refers
to as a double scree line. Catell's procedure is
modified when a double scree line is observed. Factors
are retained down to and including the factor which
begins the upper scree line. Under this modification
the correct number of factors would be retained.
Notice that Kaiser's criterion was a flawless indicator
for structure 1 and N=200. Figure 12 is another ranked
mean eigenvalue graph. This time structure 23 provides
the data. The sample size is 100. Structure 23 has 3
nuisance factors and 3 nuisance variables. Notice that
the same break occurs between the mean of the
eigenvalue magnitudes of factor numbers 3 and 4. This
time, however, the confidence intervals are quite wide
and overlap. It is not clear whether or not a break in
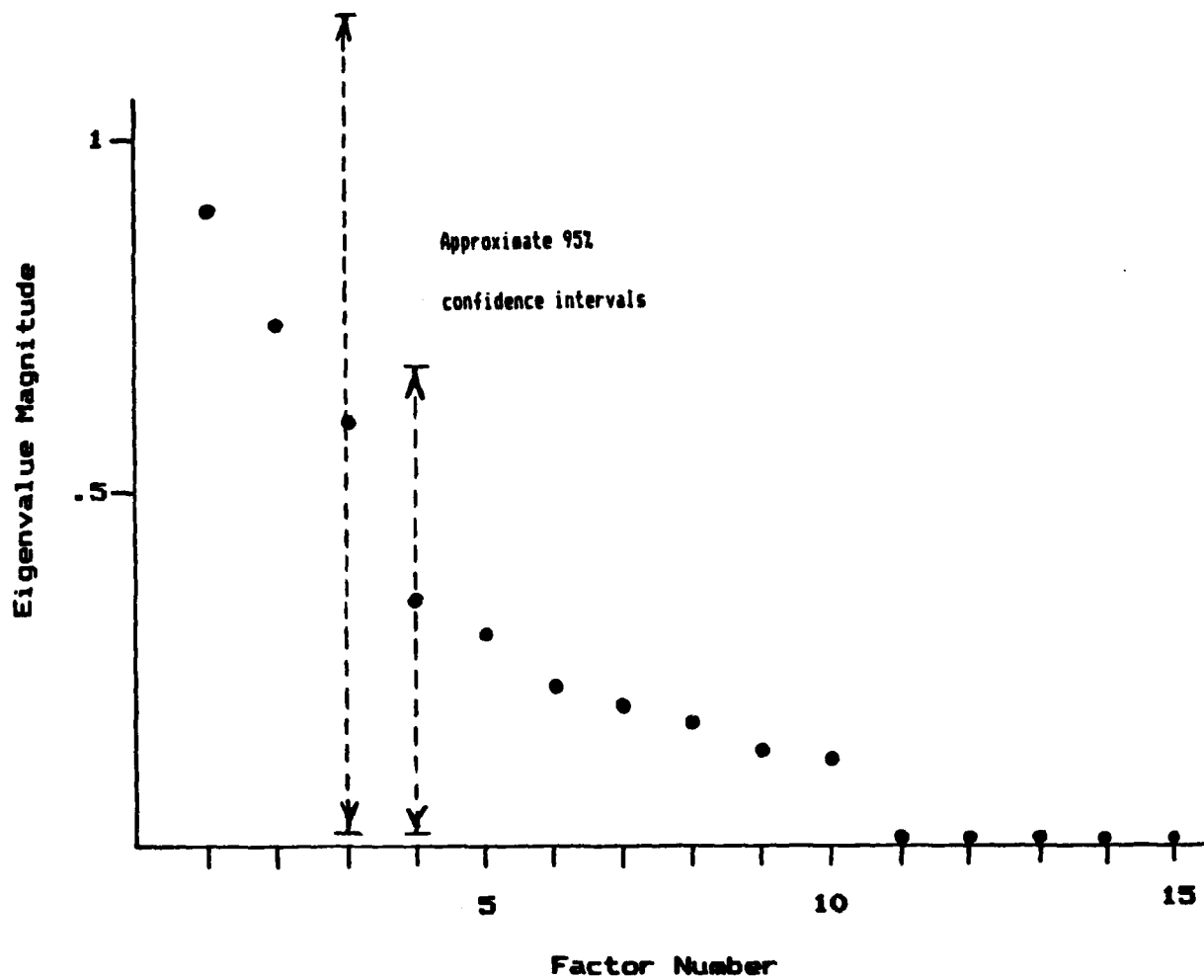the eigenvalues will even appear in a particular
sample.

Figure 12. Scree Test, Structure 23, N=100

Factor Interpretation Analysis


This section presents the results of a regression
study undertaken to determine if the sampling errors of
the experimental region were predictable.

Experimental Procedure. To reiterate the
experimental procedure, first sample vectors were
generated from the population covariance matrices.
These sample vectors were then used to form a sample
correlation matrix.  The sample correlation matrix was
factor analyzed using the PCA procedure and the
resultant factor loadings matrix (of the correct
dimensionality) was then rotated, via a least squares
procedure, to fit the original population structure.
The mean square discrepancy between the sample loadings
matrix and the population loadings matrix was then
calculated.  This mean square error was calculated
across all the loadings.  The mean square error (MSE)
is calculated by the formula

$$\frac{\sum_{j=1}^{N} \sum_{i=1}^{M} (a_{ij} - \hat{a}_{ij})^2}{M \cdot N}$$

where the $a_{ij}$ are the factor loadings for the
population factor loadings matrix, $\hat{a}_{ij}$ are the factor
loadings for the sample factor loadings matrix, M is
the number of variables, and N is the number of

55

factors. The root mean square (RMS) error is taken as
the square root of the MSE.

Performance of the Complexity Index. Figure 13 is
a plot of MSE versus the average communalities of the
original structures. Notice how structures 6 and 7
produce noticable "bumps" in the set of curves. These
two structures have four variables which load
significantly on more than one factor. All the other
structures used in figure 13 contained only univocal
variables. One would expect small bumps due to
sampling fluctuations but the aberration due to
structure 7 seems a bit severe. Figure 14 is a plot of
MSE versus the complexity index. This graph displays
more of the monotonicity one would expect. A similiar
graph is presented for the more complicated structures
which are perturbed by nuisance factors and variables.
In this graph one notices that there are two pairs of
structures whose complexity indices are quite close.
In all but one case the corresponding MSEs were quite
close. The exception occurs for structures 26 and 27.
The variance of these MSE values are of the order
.00001 and so it seems clear this particular variation
is not due to sampling error. It is probably due to
one of the complexity index's inherent weaknesses as
mentioned previously. All in all the index seems to be
performing fairly well. At least the complexity index
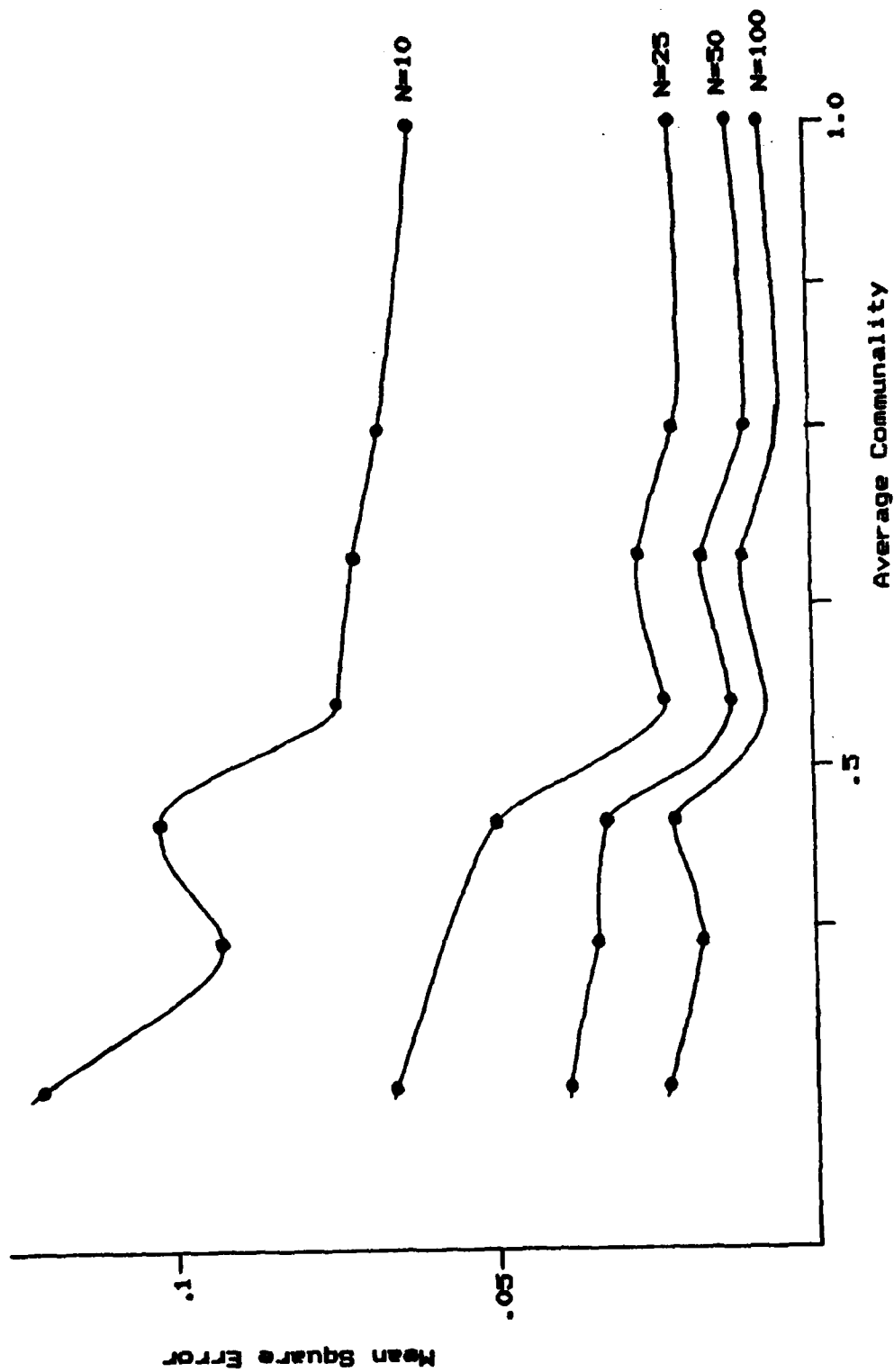is an improvement over using average communality (or

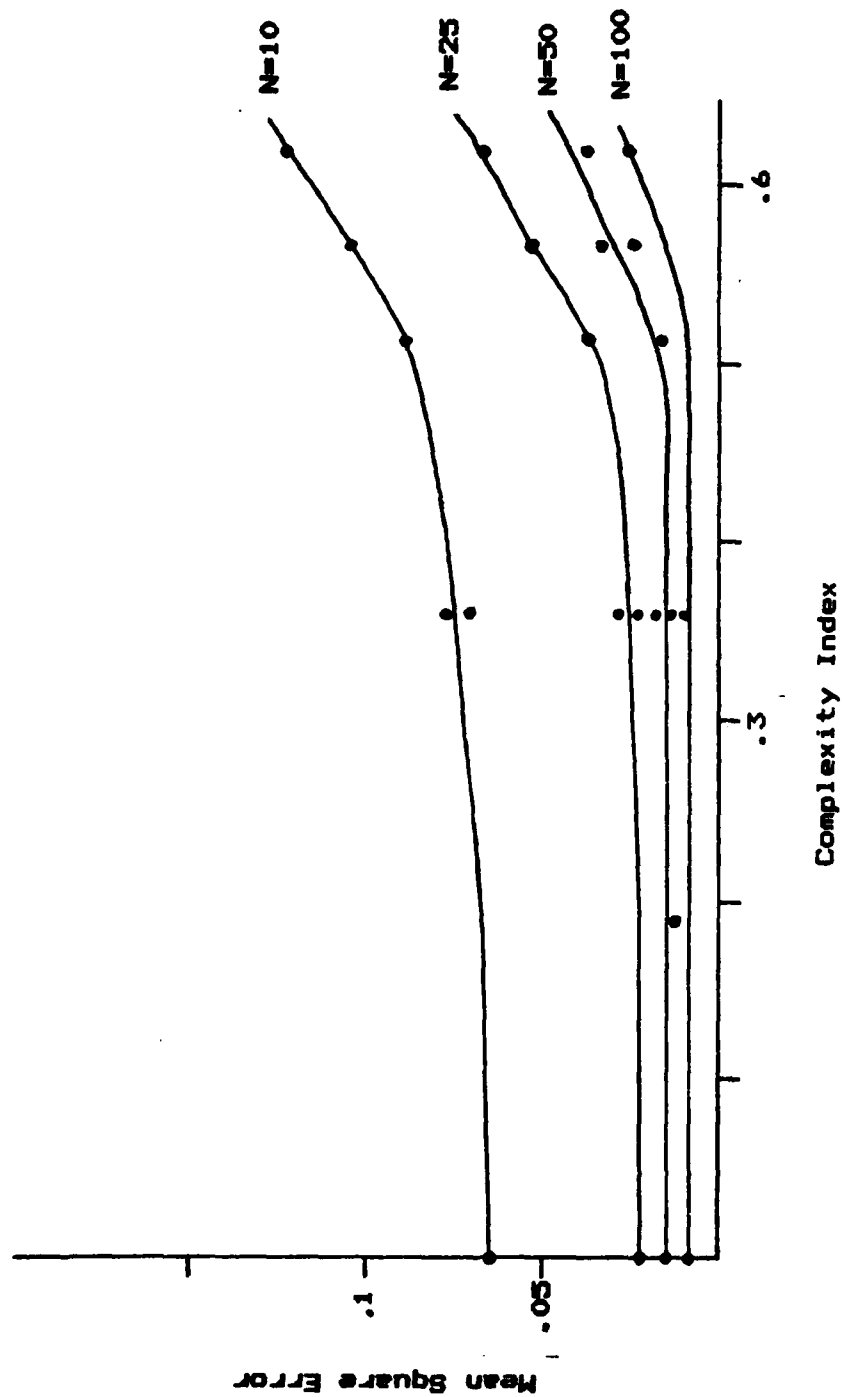56

Figure 13. Average Communality vs. Mean Square Error

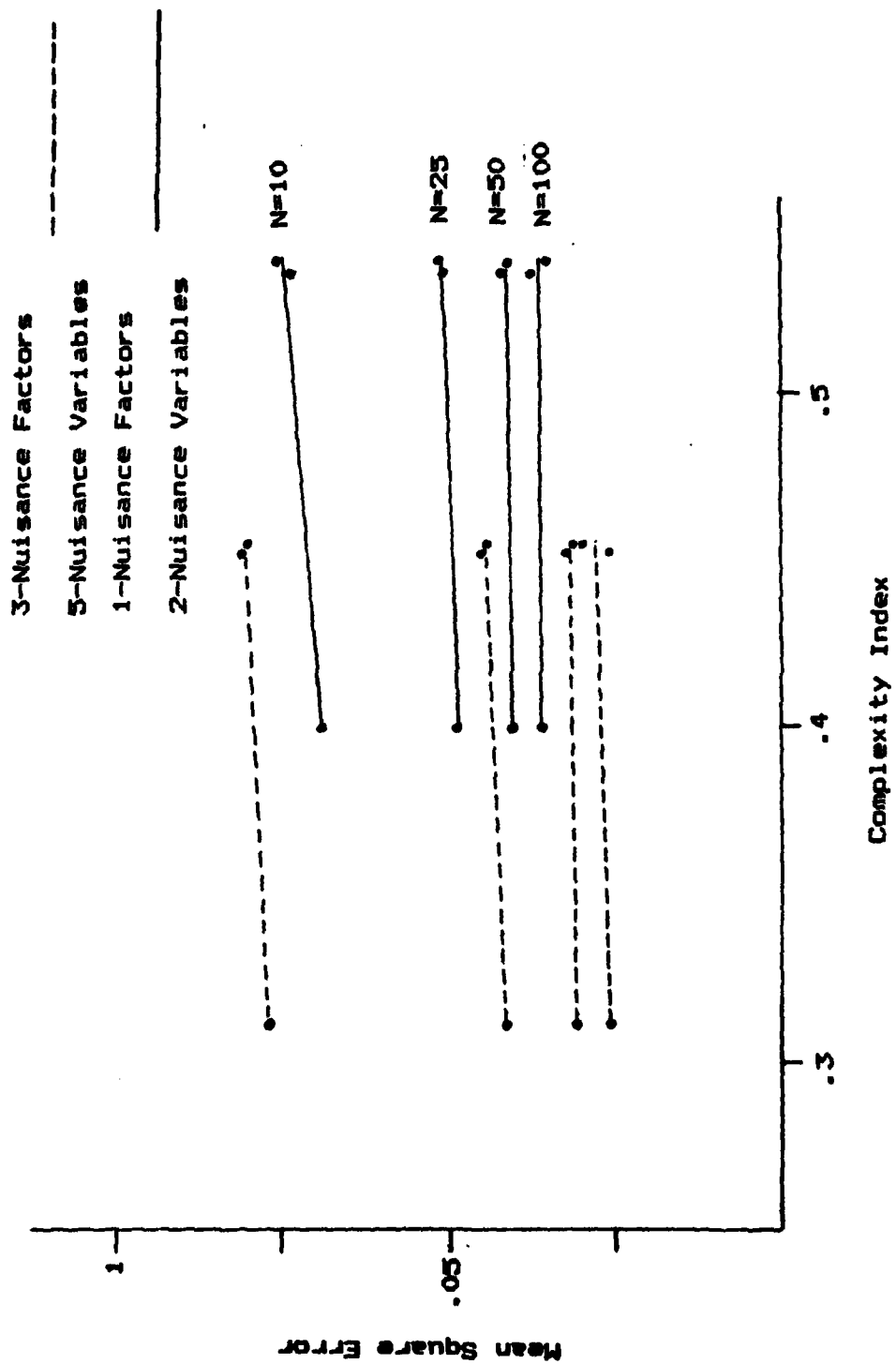Figure 14. Complexity Index vs. Mean Square Error

Figure 15. Complexity Index vs. Mean Square Error

uniqueness) as a criterion of a population structure's
complexity.

Regression Study. Several different regression
models were hypothesized and tested in order to
determine if MSE or RMS errors could be reasonably
predicted as functions of sample size, the number of
variables, the number of inherent factors, complexity
of the population structure and the interactions
between these. The condition number of the sample
correlation matrix was also examined for its possible
aid in predicting MSE or RMS errors. In these studies
each MSE or RMS value was taken as the grand mean of
1000 iterations on a particular structure-sample size
combination. The same, then, is true for each sample
condition number. Two types of regression models were
attempted.

1) Linear models with interactions--these models
were run using the Statistical Package for the Social
sciences (SPSS) (Nie, 1975). A stepwise regression
scheme was employed for variable selection. The
following models were run:
a) MSE as a linear function of sample size,
number of factors, number of variables, and all possibe
multiplicative interaction combinations. This model
was also ran with RMS as the dependent variable.
b) MSE as a linear function of sample size,
number of factors, number of variables, complexity
index of the population structure, and all possible
multiplicative interaction combinations. This model
was also ran with RMS as the dependent variable.
c) MSE as a linear function of sample size,
number of factors, number of variables, condition
number of the sample correlation matrix, and all
possible multiplicative interaction combinations. This
model was also ran with RMS as the dependent variable.

2) Nonlinear models--These models were also run
on SPSS. Nonlinear production functions of the

60

Cobb-Douglas type (Nicholson, 1978) were run over
various combinations.  The Cobb-Douglas type function
was chosen because preceived nonlinearities which were
"eyeballed" in the data.  Also, the Cobb-Douglas
function is flexible in the sense that it can mold
itself to many different shapes.  The Cobb-Douglas
function is of the form

$$Y = B_0 \ X_1^{B_1} \ X_2^{B_2} \cdots X_n^{B_n}$$

The following models were run:
        a)  RMS with independent variables:  sample
size, number of factors, and number of variables.
        b)  RMS with independent variables:  sample
size, number of factors, number of variables, and
complexity index of the population structure.
        c)  RMS with independent variables:  sample
size, number of factors, number of variables, and
condition number of the sample correlation matrix.
     The following are used as abbreviations:

     1)   Sample size--N
     2)   Number of factors--FAC or F
     3)   Number of variables--VAR or V
     4)   Complexity Index--C
     5)   Condition number--K
     6)   Interactions--an example is NxF or sample
size * number of factors

     Figure 16 is a tabular comparison of the results

from the linear models.  Note that the best predictions

are made from the model which includes the mean

condition number of the sample correlation matrices.

Note that sample size is the most significant

independent variable in all the models.  It is

interesting to note that the addition of the complexity

index into the first two models, although only

improving the model's predictability slightly, creates

a situation wherein the second most significant

independent variable is a interaction term on the

complexity index.  The coefficients of the predicitive

61

| Dependent Var | MSE | RMS | MSE | RMS | MSE | RMS |
|---|---|---|---|---|---|---|
| Independent Var | N, VAR, FAC, INTERACT | N, VAR, FAC, INTERACT | N, VAR, FAC, C, INTERACT | N, VAR, FAC, C, INTERACT | N, VAR, FAC, K, INTERACT | N, VAR, FAC, K, INTERACT |
| Multiple $R^2$ | .47069 | .59562 | .47471 | .62273 | .72168 | .76580 |
| Adjusted $R^2$ | .45966 | .58285 | .46376 | .60244 | .70671 | .75321 |
| Overall F | 42.68424 | 46.64277 | 43.37727 | 30.70097 | 48.228 | 60.81838 |
| Significance | .000 | .000 | .000 | .000 | .000 | .000 |
| Final Variables in Model F Ratio, Signif | N 80.9, .000 | N 70.6, .000 | N 81.38, .000 | N 44.72, .000 | VxK 27.9, .000 | N 62.51, .000 |
| | VAR 3.67, .058 | NxFxV 14.5, .000 | CxV 4.43, .038 | NxFxC 13.42, .000 | NxFxK 23.28, .000 | NxVxK 28.54, .000 |
| | | F 4.15, .045 | | F 10.31, .002 | N 26.06, .000 | FxK 17.29, .000 |
| | | | | FxC 6.68, .011 | FxV 4.36, .040 | NxFxK 18.93, .000 |
| | | | | NxVxC 3.51, .064 | K 3.64, .059 | NxK 5.94, .017 |
| Std. Error | .02069 | .0420 | .02062 | .041 | .01525 | .0323 |
| Std. Error / $\bar{y}$ | .48454 | .21234 | .4762 | .20728 | .35714 | .1633 |

Figure 16. Linear Models with Interactions (100 Observations)

1. MSE = (-.57129e-03 & N) + (.191474e-02 & VAR) + .45617e-01

2. RMS = (-.24277e-02 & N) + (.21258e-04 & NxFxV) + (-1.0075e-01 & FAC) + .29963

3. MSE = (-.57063e-03 & N) + (.198889e-02 & CxV) + .58935696

4. RMS = (-.20854e-02 & N) + (.24612e-04 & NxFxV) + (-.239287e-01 & FAC) + (.30027e-01 & FxC) + (-.108558e-03 & NxVxC) + .308459

5. MSE = (.428602e-04 & N) + (-.44844e-10 & NxFxK) + (-.300699e-03 & N) + (.19267e-03 & FxV) + (-.17897e-08 & K) + .38699e-01

6. RMS = (-.151499e-02 & N) + (.88877e-10 & NxVxK) + (-.761627 e-09 & FxK) + (.13306624e-04 & NxFxV) + (-.445727e-09 & NxK) + .219263

Figure 17. Linear Models - Regression Coefficients

63

equations given by the 6 models are given in figure 17.
Care should be taken when attempting to predict from
regression relationships which use the condition
number. If the sample size is less than 25, reasonable
results can not be guaranteed. The variability of the
sample condition number in the region studied was of
the order 1.0E+15 for the sample sizes of 10. In
summary, RMS errors are more accurately predicted than
MSE. Taken in pairs, the standard errors of the
estimates when normalized by their respective mean
estimates are always lower for RMS regressions than for
MSE regressions.

Figure 18. is a tabular comparison of the results
from the loglinear models. Note that these models
display slightly high adjusted r-squared values. Here
again, notice that the sample size is the most
significant independent variable. In the second model
complexity is the second most significant independent
variable. The nonlinear models are slightly superior
to their linear counterparts, the standard errors
normalized by the log of the mean estimate are in all
cases lower than the linear models.

The regression study shows that for the
experimental region studied the errors due to sampling
in a factor loadings matrix can be reasonably predicted
by either linear models or nonlinear models.

64

| Dependent Var | RMS | RMS | RMS |
|---|---|---|---|
| Independent Var | N V F | N V F C | N V F K |
| | | | |
| Multiple R | .77826 | .81536 | .78101 |
| Adjusted R | .77095 | .80715 | .77128 |
| | | | |
| Overall F | 106.464 | 99.355 | 80.245 |
| Significance | .000 | .000 | .000 |
| | | | |
| Final Variables in Model F Ratio, Signif | N 290.2, .000 | N 341.9, .000 | N 69.52, .000 |
| | V 21.79, .000 | V 10.41, .002 | V 22.27, .000 |
| | F 4.11, .046 | F 5.86, .017 | F 3.82, .054 |
| | | C 18.08, .000 | K 1.13, .291 |
| | | | |
| Std. Error | .15778 | .14478 | .15767 |
| Std. Error / ln y | -.09737 | -.08934 | -.09730 |

Figure 18. Loglinear Models
(100 Observations)

65

1. MSE = .16690 $N^{-.33018994}$ VAR$^{.44709288}$ FAC$^{.1173569}$

2. MSE = .27815 $N^{-.3289150}$ VAR$^{.30368093}$ FAC$^{.1288414}$ C$^{.1811936}$

3. MSE = .14440 $N^{-.29822542}$ VAR$^{.45224643}$ FAC$^{.1133628}$ K$^{.470337e-02}$

Figure 19. Loglinear models - Regression coefficients

## Conclusions

A limited examination of Kaiser's criterion and Catell's scree test indicates that Kaiser's criterion is usually good to within a factor and was always, for the structures and sample sizes addressed in this report, no more than 3 factors from the true dimensionality. In an ideal sampling situation Catell's scree test seems to retain one too many factors, and it is not always easy to identify the scree line.

The concept of a complexity index appears to be promising. If possible, a stronger upper bound needs to be found for the first term of the index. The possibility of unequal weights for the two terms could be investigated through some further regression studies.

The results of this research indicate that it is reasonable to estimate an overall mean error due to sampling for structures in the particular experimental region addressed by this report.

This author believes that this report has demonstrated that sample factor loadings matrices are sensitive to sample size and, more importantly, the structural complexity of a given population factor loadings structure.

The author recommends future research which would

address the sampling distribution of the complexity
index.

# Bibliography

1. Anderson, T.W. An Introduction to Multivariate Statistical Analysis. New York: John Wiley and Sons Inc., 1958.

2. Belsley, David A., et al. Regression Diagostics: Identifying Influential Data and Sourses of Collinearity. New York: John Wiley and Sons, Inc., 1980.

3. Browne, Michael W. "A Comparison of Factor Analytic Techniques," Psychometrika, Vol 33, #3, September 1968, 266-334.

4. Catell, R.B. "The Scree Test for the Number of Factors," Multivariate Behavioral Research, 1 (April, 1966), 245-276.

5. Catell, R. B., and Sullivan W. "The Scientific Nature of Factors: A Demonstration by Cups of Coffee," Behavioral Science, 1962,7,184-193.

6. Cliff, Norman. "Orthogonal Rotation to Congruence," Psychometrika, Vol. 31, #1, March 1966, 33-41.

7. Cliff, N. and Hamburger, C. D. "The Study of Sampling Errors in Factor Analysis by Means of Artifical Experiments," Psychological Bulletin, Vol 68, 1967, 430-455.

8. Cliff, N. and Pennell, R. "The Influence of Communality, Factor Strength, and Loading Size on the Sampling Characteristics of Factor Loadings," Psychometrika, Vol 32, #3, September 1967, 309-326.

9. Guttman, L. "Some Necessary Conditions for Common-Factor Analysis". Psychometrika, Vol 19, 1954, 149-161.

10. Harman, H. H. Modern Factor Analysis, 2nd ed. Chicago: University of Chicago Press, 1967.

11. Harris, R. J. A Primer of Multivariate Statistics, New York: Academic Press, 1975.

12. Horn, John L. "A Rationale and Test for the Number of Factors in Factor Analysis," Psychometrika, Vol 30, #2, June 1966, 179-185.

13. Joreskog, Karl G. Advances in Factor Analysis and Structural Equation Models, Stockholm, 1979.

14. Joreskog, K. G. Statistical Estimation in Factor Analysis. Stockholm: Almquist and Wiksell, 1963.

15. Kaiser, H. F. "The Application of Electronic Computers to Factor Analysis," Educational and Psychological Measurement, Vol 20, 1960, 141-151.

16. Lawley, D. N. and Maxwell, A. E. Factor Analysis as a Statistical Method, 2nd ed. London: Butterworth and Co., 1971.

17. Linn, Robert L. "A Monte Carlo Approach to the Number of Factors Problem," Psychometrika, Vol 33, #1, March 1968, 36-71.

18. Manners, G. E. and Brush, D. H. "A Simulation of Factor Analytic Reliability Varying Sample Size and Number of Variables," Psychological Reports, Vol 45, 1979. 471-478.

19. McNichols, Charles W. "An Introduction to Applied Multivariate Data Analysis (Course Notes), Air Force Institute of Technology, WPAFB, Ohio," 1980.

20. Naylor, Thomas, H., et al. Computer Simulation Techniques. New York: John Wiley and Sons, Inc., 1966.

21. Nicholson, Walter. Microeconomic Theory. Hinsdale, Illinois: Dryden Press, 1978.

22. Nie, N. H., et al. Statistical Package For The Social Sciences, 2nd ed. New York: McGraw Hill, 1975.

23. Odell, P. L., and Feiveson, G. H. "A Numerical Procedure to Generate a Sample Covariance Matrix," American Statistical Association Journal, March 1966, 199-203.

24. Pennell, Roger. "The Influence of Communality and N on the Sampling Distributions of Factor Loading," Psychometrika, Vol 33, #4, december 1968, 423-439.

25. Saunders, D. R. "Further Implications of Mundy-Castles Correlations between EEGand Weshsler-Bellevue Variables." Journal of the National Institute for Personnel Research (Johannesborg), Vol 8, 91-101.

26. Schoneman, Peter H. "A Generalized Solution of the Orthogonal Procrustes Problem," Psychometrika, Vol 31, #1, March 1966, 1-11.

27. Tucker, Ledyard R. "Evaluation of Factor Analytic Research Procedures by Means of Simulated Correlation Matrices," Psychometrika, Vol 34, #4, December 1969, 420-459.

28. Tucker, L. R. "Recovery of Factors from Simulated Data." Paper presented at the meeting of the Psychometric and Psychomomic Societies, Niagara Falls, Ontario, October 1964.

29. Westlake, Joan R. A Handbook of Numerical Matrix Inversion and Solution of Linear Equations. New York: John Wiley and Sons, Inc., 1968.